

Input Decimated Ensembles

Kagan Tumer and Nikunj C. Oza
NASA Ames Research Center
Moffett Field, CA 94035
{kagan, oza}@ptolemy.arc.nasa.gov

November 27, 2001

Abstract

Using an ensemble of classifiers instead of a single classifier has been shown to improve generalization performance in many pattern recognition problems. However, the extent of such improvement depends greatly on the amount of correlation among the errors of the base classifiers. Therefore, reducing those correlations while keeping the classifiers' performance levels high is an important area of research. In this article, we explore **input decimation (ID)**, a method which selects feature subsets for their ability to discriminate among the classes and uses them to decouple the base classifiers. We provide a summary of the theoretical benefits of correlation reduction, along with results of our method on two underwater sonar data sets, three benchmarks from the Proben1/UCI repositories, and two synthetic data sets. The results indicate that input decimated ensembles (IDEs) outperform ensembles whose base classifiers use all the input features; randomly selected subsets of features; and features created using principal components analysis, on a wide range of domains.

1 Introduction

Using an ensemble of classifiers instead of a single classifier has been repeatedly shown to improve generalization performance in many pattern recognition problems [9, 17, 69]. It is well-known that, in order to obtain such improvement, one needs to simultaneously maintain a reasonable level of performance in the base classifiers that constitute the ensemble and reduce their correlations [1, 29, 44, 49, 64]. There are many ensemble methods that actively promote diversity (e.g., modify error surface, lower correlations in the outputs) among their base classifiers [49, 58, 64]. Most work in this field, however, focuses on pattern-level selection (e.g., Bagging [9], Boosting [22]). These methods bring about diversity in the base models by training them with different subsets of the training set. One drawback of such methods is that by definition, only a portion of the available data is used during learning. This can lead to poor performance, particularly when the data sets are small to begin with. Training

the base classifiers using different subsets of features avoids this issue as all the patterns can be used in the training while still yielding base model diversity.

Two possible feature selection/extraction methods are Principal Component Analysis (PCA) [31, 52] and random subspace selection [10, 25]. PCA constructs new features such that the data has maximum variability over those features. However, PCA, when used in combining, not only generates the same features for all potential classifiers in the pool, but also fails to take class information into account. Random subspace selection overcomes the first shortcoming of PCA, but it too does not consider the class labels when generating the feature subsets. These two methods do not attempt to choose features in a manner that is helpful in the classification task.

In this paper, we present input decimation—a method of choosing different subsets of the original features based on the correlations between individual features and class labels, and training classifiers on those subsets prior to combining. This method not only reduces the dimensionality of the data, but uses this dimensionality reduction to reduce the correlations among the classifiers in an ensemble, thereby improving the classification performance of the ensemble [51, 63, 67].

Our results indicate that input decimation reduces the error up to 90% over single classifiers and ensembles trained on all features, randomly-selected subsets of features, and principal components. While we expected strong ensemble performance, input decimation also provided improvements in the base classifiers in many cases by pruning extraneous or irrelevant features, thus simplifying the learning problem faced by each base classifier. In this study we focus on the “averaging” combiner for two reasons: (i) despite its simplicity (or perhaps because of it) this combiner has been shown to perform well and hold its own against a wide array of more sophisticated methods [16, 17]; and (ii) by choosing a simple combiner we isolate the effects of input decimation from those of the combining method. Furthermore, pattern-level ensemble methods such as bagging, boosting, and stacking can be used *in conjunction* with input decimation which is a feature-level ensemble method (i.e., input decimation is orthogonal to those methods). Therefore, one can make meaningful comparisons between averaging combiners *with* and *without* input decimation, or say, between stacking or bagging *with* and *without* input decimation (not reported in this article), not between input decimated ensembles and bagging or boosting.

In Section 2, we summarize a theory of classifier ensembles that highlights the connection between correlation among base classifiers and ensemble performance, along with a brief overview of different dimensionality reduction methods. In Section 3 we present the details of the input decimated ensemble, and in Section 4 we provide experimental results on two real underwater sonar data sets, three data sets from the PROBEN1/UCI benchmarks [6, 54], and two synthetic data sets which allow a systematic study of input decimation. We conclude with a discussion on the effectiveness of input decimation under various circumstances along with future research directions in Section 5.

2 Background

Model selection is a ubiquitous problem in many pattern recognition problems. Neither the selection of the method (e.g., multi-layer perceptron, nearest neighbor algorithm), nor the tuning of that algorithm can yet be fully automated for all problems [15, 20, 23]. The use of ensembles provides partial relief since by pooling the classifiers before a decision is made, potential sensitivity to any single model is greatly reduced. Of course, the more similar the classifiers are, the less likely it is that new information will be present in the ensemble, resulting in little more than a “rubber stamping” committee. In this section we first formalize this connection between the correlation among the classifiers’ errors and ensemble performance and then discuss various methods that aim to reduce that correlation.

2.1 Correlation and Ensemble Performance

In this article we focus on classifiers that model the *a posteriori* probabilities of the output classes. Such algorithms include Bayesian methods [4], and properly trained feed forward neural networks such as Multi-Layer Perceptrons (MLPs) [56]. We can model the i th output of such a classifier as follows (details of this derivation are in [63, 64]):

$$f_i(x) = P(C_i|x) + \eta_i(x),$$

where $P(C_i|x)$ is the posterior probability of the i th class given pattern x , and $\eta_i(x)$ is the error associated with the i th output. Given an input x , if we have one classifier, we classify x as being in the class i whose value $f_i(x)$ is largest.

Instead, if we use an ensemble that calculates the arithmetic average over the outputs of N classifiers $f_i^m(x)$, $m \in \{1, \dots, N\}$, then $P(C_i|x)$ is given by:

$$f_i^{ave}(x) = \frac{1}{N} \sum_{m=1}^N f_i^m(x) = P(C_i|x) + \bar{\eta}_i(x), \quad (1)$$

where:

$$\bar{\eta}_i(x) = \frac{1}{N} \sum_{m=1}^N \eta_i^m(x)$$

and $\eta_i^m(x)$ is the error associated with the i th output of the m th classifier.

Now, the variance of $\bar{\eta}_i(x)$ is given by [64]:

$$\begin{aligned} \sigma_{\bar{\eta}_i}^2 &= \frac{1}{N^2} \sum_{l=1}^N \sum_{m=1}^N \text{cov}(\eta_i^l(x), \eta_i^m(x)) \\ &= \frac{1}{N^2} \sum_{m=1}^N \sigma_{\eta_i^m}^2(x) + \frac{1}{N^2} \sum_{m=1}^N \sum_{l \neq m}^N \text{cov}(\eta_i^l(x), \eta_i^m(x)). \end{aligned}$$

If we express the covariances in terms of the correlations ($cov(x, y) = corr(x, y)\sigma_x\sigma_y$), assume the same variance $\sigma_{\eta_i}^2$ across classifiers, and use the average correlation factor among classifiers, δ_i , given by

$$\delta_i = \frac{1}{N(N-1)} \sum_{m=1}^N \sum_{l \neq m} corr(\eta_i^l(x), \eta_i^m(x)), \quad (2)$$

then the variance becomes:

$$\sigma_{\eta_i}^2 = \frac{1}{N} \sigma_{\eta_i(x)}^2 + \frac{N-1}{N} \delta_i \sigma_{\eta_i(x)}^2 = \frac{1 + \delta_i(N-1)}{N} \sigma_{\eta_i(x)}^2. \quad (3)$$

Based on this variance, we can compute the variance of the decision boundary and, generalizing this result to the classifier error, we obtain the relationship between the model error (beyond the Bayes error) of the ensemble (E_{model}^{ave}) and that of an individual classifier (E_{model}) [63, 64]:

$$E_{model}^{ave} = \left(\frac{1 + \delta(N-1)}{N} \right) E_{model} \quad (4)$$

where

$$\delta = \sum_{i=1}^L P_i \delta_i \quad (5)$$

and P_i is the prior probability of class i .

Equation 4 quantifies the connection between error reduction and the correlation among the errors of the base classifiers. This result leads us to seek to reduce the correlation among classifiers prior to using them in an ensemble.

2.2 Correlation Reduction Methods

As shown above, if the classifiers to be combined repeatedly provide the same (either erroneous or correct) classification decisions, there is little to be gained from combining, regardless of the chosen scheme. As equation 4 shows, reducing δ and increasing N are two ways to improve the performance of a classifier ensemble. However, these two ways are not independent. This phenomenon is best illustrated by Figure 1, where the error reduction depending on the correlation among the classifiers is displayed as a function of the number of classifiers (based on Equation 4). For example, even though increasing the number of classifiers from 4 to 8 does not provide any sizeable gains when the correlation is .9, it provides significant gains if the correlation is .1. That is, keeping the correlations low not only provides better error reduction for a given number of classifiers, but provides greater gains when adding classifiers.

To improve ensemble performance one must either actively promote diversity during training or achieve diversity through the selection of the data presented to

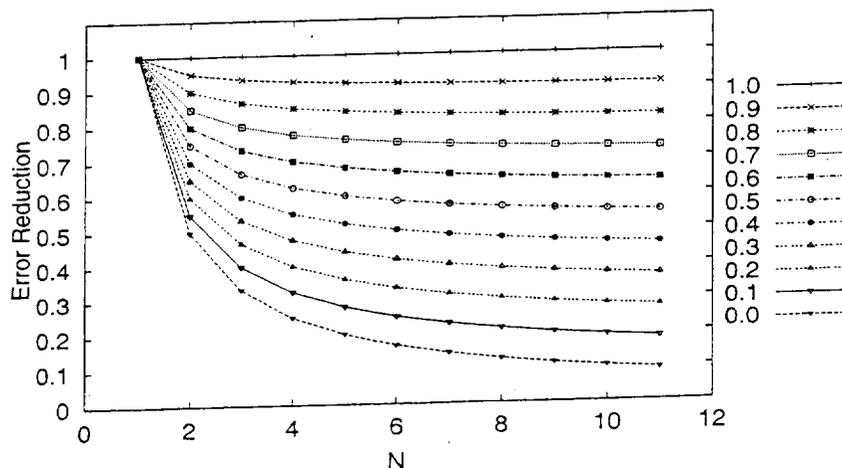


Figure 1: Effect of correlation on error reduction.

the base classifier training algorithms. Examples of the former include distorting the output space through error-correcting output codes [18], using principal component analysis on the output space [45], using genetic algorithms to train the classifier [50, 61] or modifying the error function used for training [58]. Examples of the latter include bagging [9], cross-validation partitioning [40, 64] and even boosting [22] (though the goal there is not to reduce correlation, the net effect is the same). The most common data selection methods focus on the “pattern” space, though dimensionality reduction methods which manipulate the feature space can also be used. Feature space methods have the advantage that they do not reduce the number of patterns available for training each classifier. They generally fall into one of two different classes of methods: feature selection or feature extraction.

Feature extraction algorithms such as Principal Components Analysis (PCA) [5, 31, 52] or Independent Component Analysis (ICA) [28] reduce the dimensionality of the data by creating new features. Linear PCA, perhaps the most commonly used feature extraction method, creates new features that are linear combinations of the original features. The aim of PCA, however, is to devise features on which the data shows the highest variability, whether those features are useful for classification or not [5]. Furthermore, because all the information present in the initial features is “crammed” into fewer principal components, there is a danger that classifiers trained on the principal components will have higher, not lower correlations among them. Figure 2 demonstrates the perils of not using class information. The left half of the figure shows a case in which PCA works effectively. In this case the first principal component (y_1) corresponds to the variable with the highest discriminating power. The right half shows a similar data set (similar data distribution and linearly separable). However, because the first principal component is not “aligned” with the class labels, selecting this

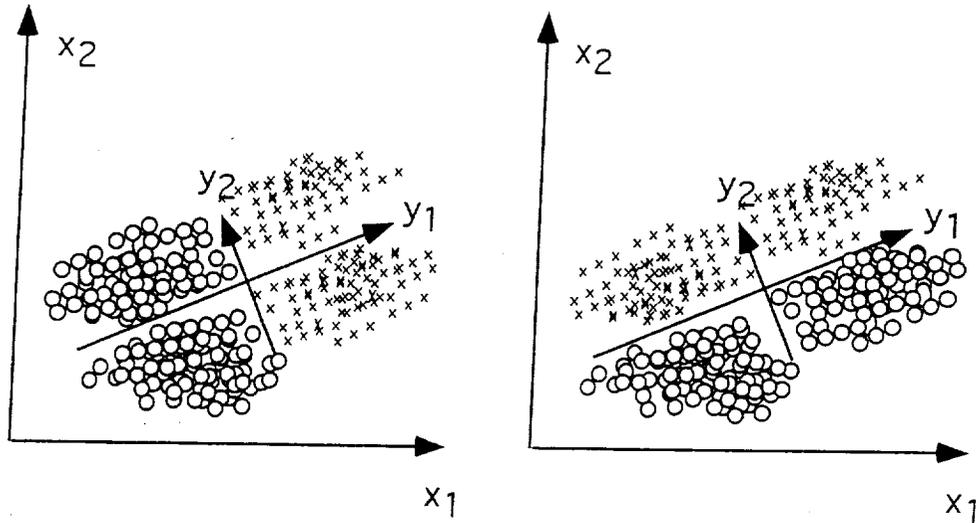


Figure 2: PCA and classification: The first principal (y_1) can provide a good discriminating feature (left) or a poor one (right), since the class membership information is not used.

component is a poor choice for this problem. Indeed, an input set consisting of only the first component would provide practically random decisions on this data set.¹ Yet, PCA remains one of the most frequently used dimensionality reduction methods in many classification domains, including medical and space applications [55, 60].

Feature selection algorithms focus on selecting a subset of the features to present to the classifiers. One example is the random subspace method [10, 27] where random subsets of the original features are presented to the classifiers. However, looking at y_1 and y_2 (assuming those two are the original features) in Figure 2 shows a pitfall of random feature selection. Randomly selecting feature y_1 in the class configuration shown in (a) will lead to satisfactory classification, whereas randomly selecting feature y_1 in (b) will lead to all discriminating information being lost. Many other feature selection methods use various criteria for deciding the relevance of each feature to the task at hand and choose some subset of the features according to those criteria [3, 7, 8, 19, 30, 43]. The subset selection can be distinct from the learning, which is the case with *filter* methods. However, most of these feature selection methods attempt to choose features that are useful in discriminating across all classes. Using such a method within an ensemble learning scheme would have limited effectiveness since it would choose the same features for every base classifier, leading to rela-

¹There are variations on PCA that use local and/or nonlinear processing to improve dimensionality reduction [13, 33, 34, 47, 48, 59]. Although these methods implicitly account for some class information and therefore are better suited than global PCA methods for classification problems, they do not directly use class information.

tively small correlation reduction. One exception is to break an L -class problem into $\binom{L}{2}$ two-class problems and perform feature selection within each of those problems [41]. In many real-world problems, there are features that are useful at distinguishing whether a pattern is of one particular class but are not useful at distinguishing among the remaining classes. In the next section we present input decimation, which takes advantage of this fact to reduce both the dimensionality and correlation in classifier ensembles.

3 Input Decimated Ensembles

Input Decimation decouples the classifiers by exposing them to different aspects of the same data. ID trains L classifiers, one corresponding to each class in an L -class problem.² For each classifier, the method selects a user-determined number of the input features having the highest absolute correlation to the presence or absence of the corresponding class.³ The objective is to “weed” out input features that do not carry strong discriminating information for a particular class, and thereby reduce the dimensionality of the feature space to facilitate the learning process. Additionally, the classifiers’ features are selected using different relevance criteria, which leads to different feature subsets for each base classifier and a reduction in their correlations.

Let the training set take the following form:

$$\{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_m, \mathbf{y}_m)\},$$

where m is the number of training examples. Each \mathbf{x}_i has $\|F\|$ elements (where F is the set of input features) representing the values of the input features in example i . Each \mathbf{y}_i represents the class using a distributed encoding, i.e., it has L elements, where L is the number of classes, $y_{il} = 1$ if example i belongs to class l and $y_{il} = 0$ otherwise. In this study our base classifiers consist of MLPs trained with the backpropagation algorithm.⁴

Given such a data set, and a base classifier learning algorithm, input decimated ensembles operate as follows:

- For each class $l \in \{1, 2, \dots, L\}$,
 1. Compute the absolute value of the correlation between each feature j (x_{ij} for all patterns i) and the output for class l (y_{il} for all patterns i).
 2. Select the n_l features having the highest absolute correlation, resulting in new feature set F_l . One can either predetermine n_l based on

²More generally, one trains nL classifiers where n is a positive integer.

³Note that this method requires the problem to have at least three classes. In a two-class problem, features strongly correlated with one class will be strongly anti-correlated with the other class, so the same features would be chosen for both classifiers.

⁴In principle, any learning algorithm that estimates the a posteriori class probabilities can be used.

prior information about the data set, or learn the value to optimize performance.

3. Construct a new training set by retaining only those elements of the x_i 's corresponding to the features F_l and all the outputs.
4. Run the base classifier learning algorithm on this new training set. Call the resulting classifier f^l .⁵

Given a new example x , we classify it as follows:

- For each class $k \in \{1, 2, \dots, L\}$, calculate $f_k^{ave}(x) = \frac{1}{L} \sum_{l=1}^L f_k^l(x)$, by presenting the proper features F_l of example x to each of the L classifiers f^l .⁶
- Return the class $K = \operatorname{argmax}_k f_k^{ave}(x)$.

Fundamentally, input decimation seeks to reduce the correlations among individual classifiers by using different subsets of input features, while pattern-level methods such as bagging and boosting attempt to do so by choosing different subsets of training patterns.

4 Experimental Results

In this section, we present the results of input decimation on two underwater sonar data sets, three Proben1/UCI benchmark data sets and two synthetic data sets. In all results reported below, the base classifiers consist of Multi-Layer Perceptrons (MLPs) with a single hidden layer trained with the backpropagation algorithm. The learning rate, momentum term, and number of hidden units were experimentally determined. In all cases, we report test set error rates averaged over 20 runs, along with the differences in the means.⁷

4.1 Passive Sonar Signals

A real world problem with all the characteristics required for a complete study is that of classifying short duration underwater signals obtained from passive sonar signals [14]. Both biological and non-biological phenomena produce such short duration sounds, and experts can determine the cause by studying their pulse signatures or spectrograms. Automating this classification process is a difficult process because these signals are highly non-stationary, have different spectral characteristics depending on sources or propagation paths and may

⁵If one is training nL classifiers for $n > 1$, then the algorithm calls the base classifier learning algorithm n times to create n classifiers $f^{l1}, f^{l2}, \dots, f^{ln}$ with feature set F_l .

⁶If we are instead training nL classifiers for $n > 1$, then we calculate $f_k^{ave}(x) = \frac{1}{nL} \sum_{l=1}^L \sum_{i=1}^n f_k^{li}(x)$.

⁷That is, for an error with mean μ and variance σ^2 , we report the $\mu \pm \sigma/\sqrt{K}$ where K is the number of repetitions ($K=20$ for experiments reported here). Confidence intervals of desired sensitivity can be obtained directly from the differences in the means.

have significant overlap. A more detailed description of the sonar signals and the difficulty associated with their classification can be found in [24, 63].

The two data sets used for this experiment are both extracted from short-duration passive sonar signals due to four naturally occurring oceanic sources (sound of ice cracking, porpoise and two different whale sounds). Although there is some complementarity among the data sets, for the purposes of this study we will treat them as different data sets.⁸ The first set, SONAR1, consists of 25 features, including 16 Gabor wavelet coefficients,⁹ signal duration and other temporal descriptors and spectral measurements. There were 496 training and 823 test patterns. The second set, SONAR2, consists of 24 features, including reflection coefficients corresponding to the maximum broadband energy segment using both short and long time windows, signal duration and other temporal descriptors. There were 564 training and 823 test patterns. For both data sets, we used an MLP with 50 hidden units.

Tables 1 shows the error rates, differences in the mean, and correlation among the base classifiers for both the full feature set and the input decimated set. In this case, each base classifier had an input decimated set of 22 features for both SONAR1 and SONAR2 after features with little correlation to each output were deleted. Retaining more features did not result in a significant drop in correlations, whereas removing more features resulted in drops in individual classifier performance that were too large to be compensated by combining. In fact, this data set is not particularly well-suited for input decimation because it has a small number of carefully-extracted, relevant features.

Table 1: Ensemble Performance on both sonar data.

		Full Feature Set		Input Decimation	
		Error Rate	δ	Error Rate	δ
SONAR1	1	7.47 \pm .10	.89	8.38 \pm .15	.68
	4	7.05 \pm .07		7.10 \pm .07	
	8	7.17 \pm .05		6.99 \pm .06	
SONAR2	1	9.95 \pm .16	.76	9.73 \pm .16	.72
	4	9.26 \pm .15		8.80 \pm .06	
	8	8.94 \pm .11		8.62 \pm .06	

For SONAR1, the deletion of even lowly-correlated inputs affects the performance of the base classifier significantly. However, due to the correspondingly large reduction in the error correlation, input decimated ensembles perform at the level of the full feature set, with IDE for $N = 8$ providing a statistically significant gain over the full feature set ensemble and IDE for $N = 4$ at the $\alpha = .05$ level. For SONAR2, the gains are more significant in that even the input decimated base classifier improves slightly upon the full featured base classifier, allowing for sizable gains by the input decimated ensemble. This is

⁸See [63] for a study where the two data sets were used in conjunction.

⁹Gabor wavelet coefficients provide a multiscale representation that does not assume signal stationarity [12].

achieved in spite of the relatively modest drop in the error correlation among the base classifiers. Also, note that for SONAR1, because the correlation is high for the base classifiers trained on the full feature set, increasing the number of classifiers from 4 to 8 does not provide any gains (instead it provides statistically equivalent errors).

4.2 Proben1/UCI Benchmarks

In the SONAR data presented above each feature carried a significant amount of discriminating information. In fact, because each feature was carefully extracted from the raw data, one should not have expected much improvement through input decimation. In this section we perform a more detailed analysis on three benchmark data sets where we gradually decrease the dimensionality until we end up with 5-10% of the original features. On these benchmark sets, we expect this more extreme case of input decimation to expose the strengths and weaknesses of this method.

The three data sets from the UCI/PROBEN1 benchmarks [6, 54] selected for this study were: The Gene dataset from the PROBEN1 (i.e., using train/test split from PROBEN1), and the Splice junction gene sequences and Satellite Image datasets (Statlog version) from the UCI Machine Learning Repository. The Gene data set has 120 input features and three classes [46, 54]. The MLP has a single hidden layer of 20 units, a learning rate of 0.2 and a momentum term of 0.8. The Splice data consists of 60 input features and three classes [6]. Here we selected an MLP with a single hidden layer composed of 120 units, a learning rate of 0.05, and a momentum term of 0.1. The Satellite Image data set has 36 input features and 6 classes [6]. We selected an MLP with a single hidden layer of 50 units, and a learning rate and momentum term of 0.5. The ensembles consisted of three classifiers for Gene and Splice and six classifiers for Satellite Image—the same as the number of classes.

Figures 3-5 show the classification performance and classifier correlations for all three data sets, averaged over 20 runs. For clarity we omit the error bars, since they ranged from 0.05 to 0.25% and as such were smaller than the symbols representing the data points. The rightmost point in each graph (e.g., the point corresponding to 120 features for the Gene data set) shows the full feature set performance. For the Gene data, the full feature ensemble is significantly more accurate than the single classifier, while for the Satellite Image and Splice data sets, the ensemble is only marginally more accurate.

In case of the Gene data, the average ensembles with 20, 30, and 40 inputs are significantly more accurate than both the original network ensembles described in the previous section and their PCA counterparts. With IDE, the performance of the ensemble goes up as the number of features increases until all the relevant features are included and then starts declining with the addition of irrelevant features. The average correlation behaves the same way. For 10 or fewer features, we expect the average correlation to be low because different sets of 10 features have the highest relevance to each class. As the number of features increases up to 30, the base classifiers have an increasing number of common

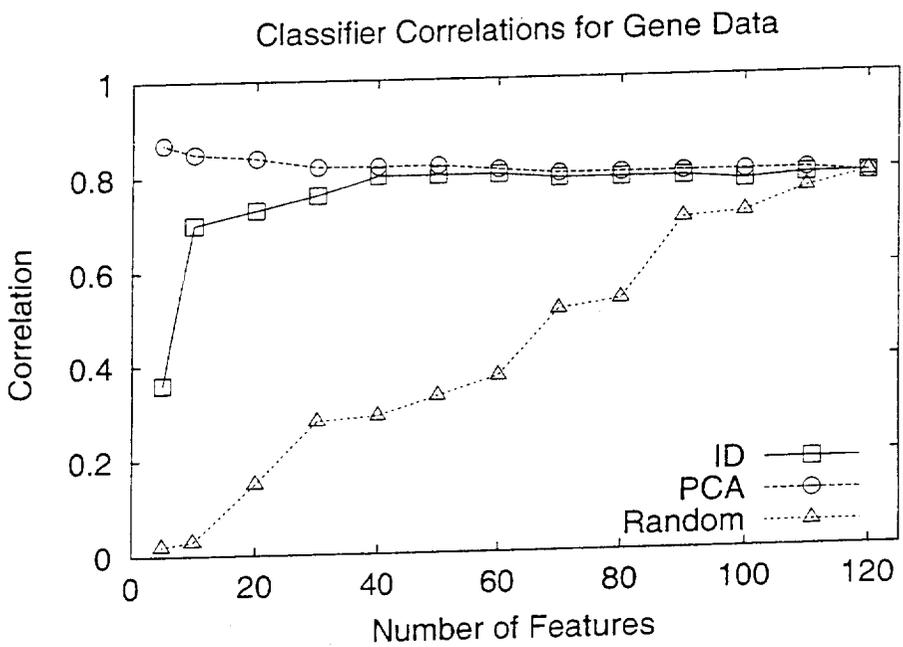
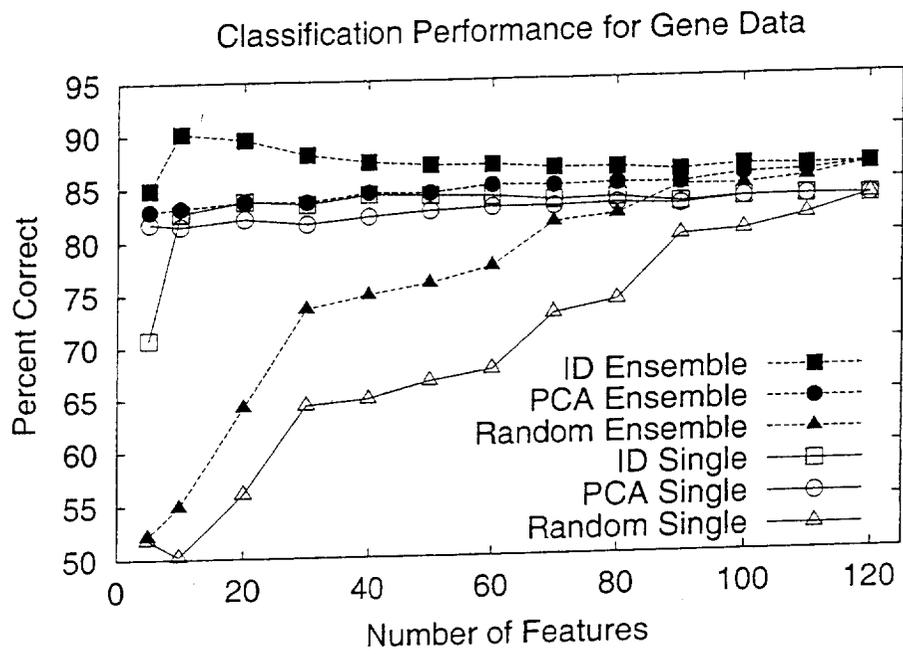


Figure 3: Performance (a) and Correlations (b) for the Gene data set.

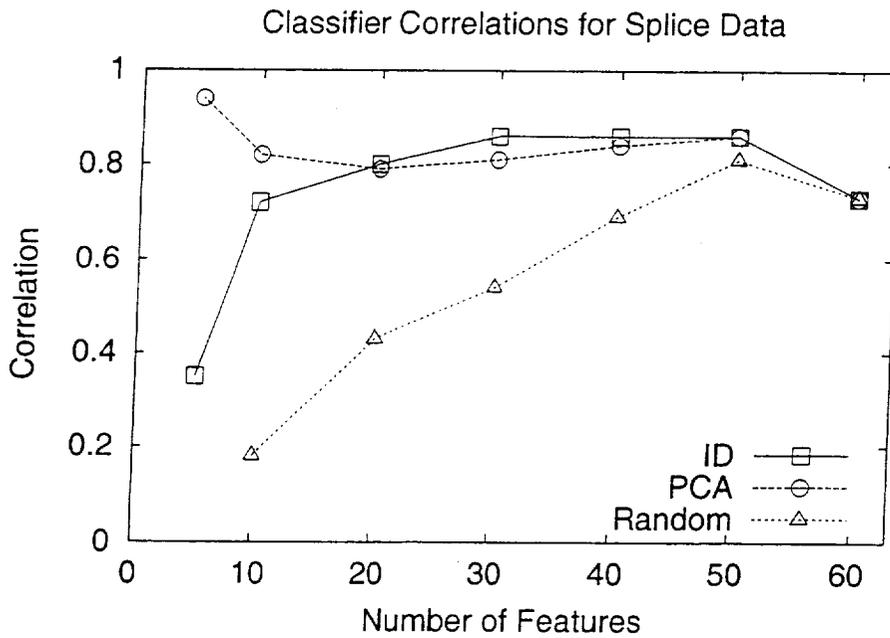
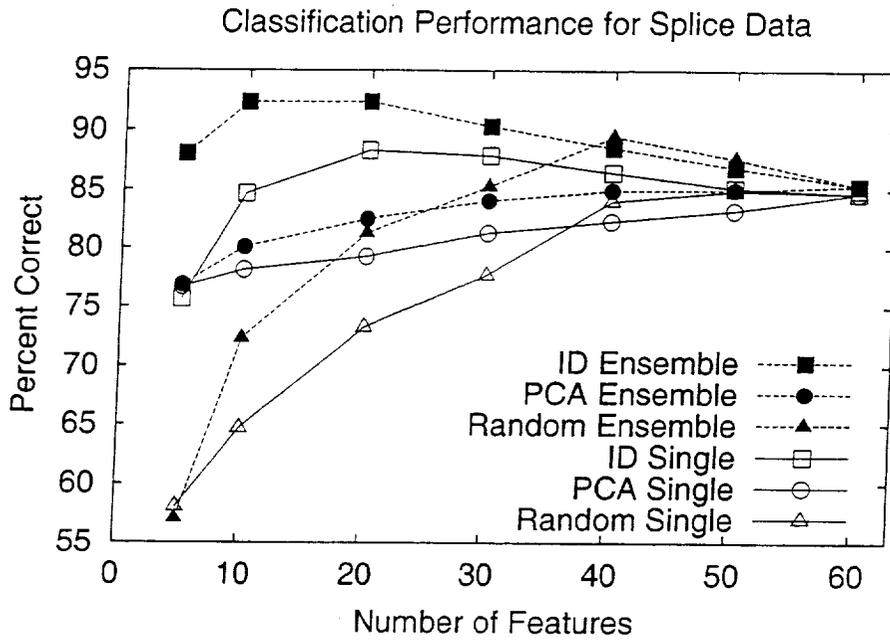


Figure 4: Performance (a) and Correlations (b) for the Splice data set.

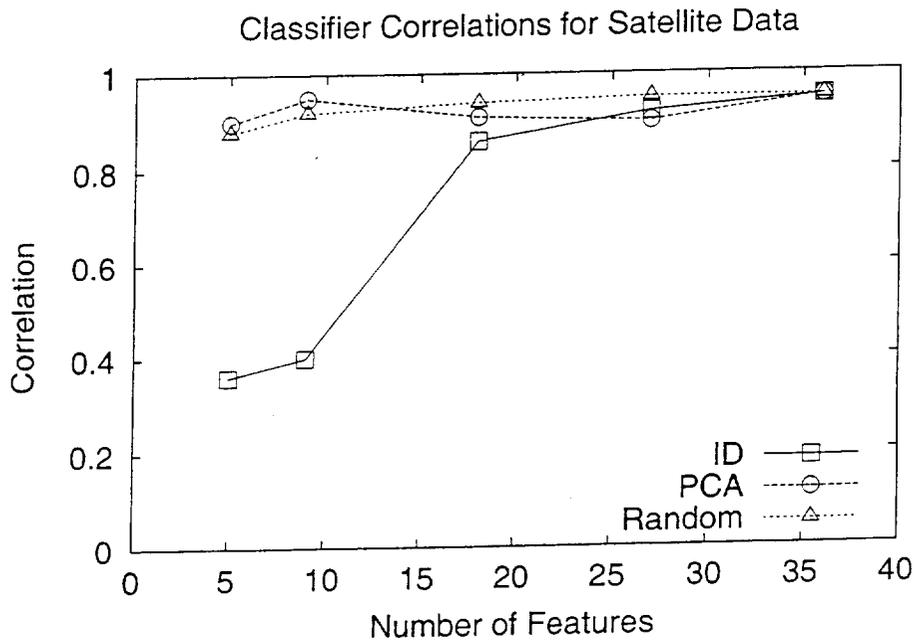
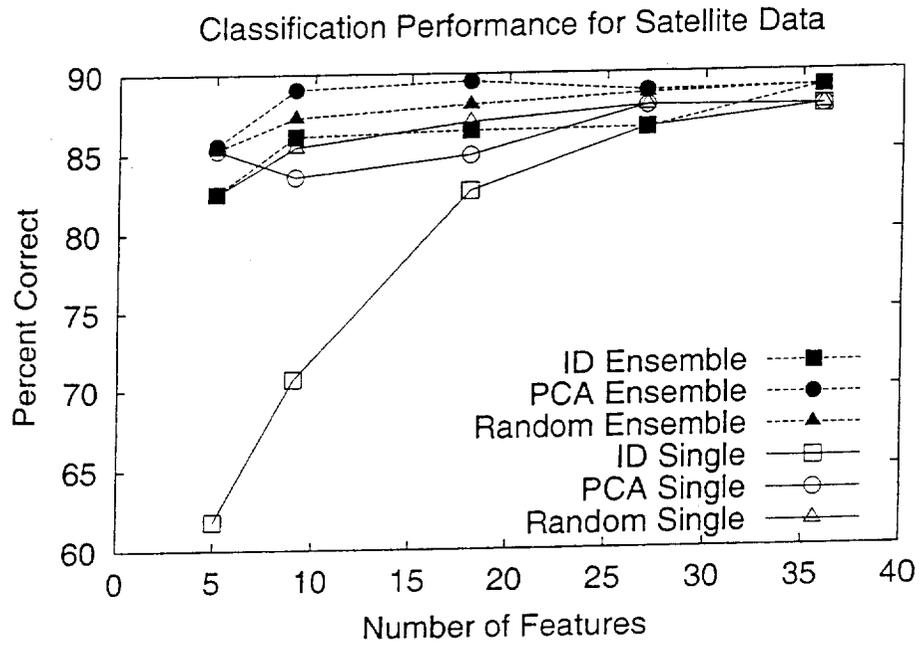


Figure 5: Performance (a) and Correlations (b) for the Satellite data set.

features. At 30, the base classifiers have virtually all common features—the 30 that are relevant to all three classes; therefore, we would expect maximum average correlation. Beyond 30, each base classifier is getting (probably different) irrelevant features, leading to a reduction in correlation. With PCA, the performance of the ensemble is relatively stable and inferior to IDE. This is consistent with the fact that principal components are not necessarily good discriminative features, and adding principal components beyond the first few would likely have little effect on the classification performance. The performance of the ensemble with random feature subsets increases in random increments with the addition of features depending on how relevant they are. On this dataset, the performance of random feature ensembles was uncompetitive because random selection never yielded good feature subsets.

In the Splice data experiments, *all* the decimated feature-based ensembles significantly outperformed both the original ensemble and the PCA-based ensembles. Random feature-based ensembles performed somewhat better here than in the Gene data set. With 40 and more features, it was competitive to input decimation. However, the best performing predictor overall is clearly the input-decimated ensemble with 10 inputs per classifier. What is particularly notable in this case is that a reduction of dimensionality based on PCA has a strong negative impact on the classification performance. With 20 principal components for example, the performance of the single classifiers drops by 7% relative to the single classifier with all the input features, whereas the performance of the ID single classifier increases by 3%. The improvement of the performance of the single classifiers due to decimation is an initially surprising aspect of these experiments since one may not expect to find too many “irrelevant” features in these real data sets. However, an analysis shows that the inputs that were decimated were in fact providing “noise” to the classifier. Although it is theoretically true that the classifier with more information will do at least as well as the classifier with less information, in practice with only a limited amount of data, extracting the correct information can cause a problem for such classifiers causing them to perform worse than their counterparts with less information.

On the Satellite Image data however, the input decimated ensemble with 27 features was the only one that did not perform significantly worse than the single classifier and the original ensemble. Both the PCA and random feature ensembles outperformed IDE. Because the single IDE classifiers performed much worse than the PCA and random feature single classifiers, we examined the features that were chosen in each ensemble. Figure 6 shows the average correlations among the features chosen for the base classifiers in the three types of ensembles. The features that IDE chose have a much higher correlation among themselves relative to random and PCA ensembles, especially for smaller numbers of inputs. This means that IDE often chooses several features with high correlations to the class without realizing that they may be redundant. Random feature selection does not fall into this trap since it does not consider correlations at all. PCA’s correlations are the lowest because it creates features specifically designed to have low correlations among each other. Among the

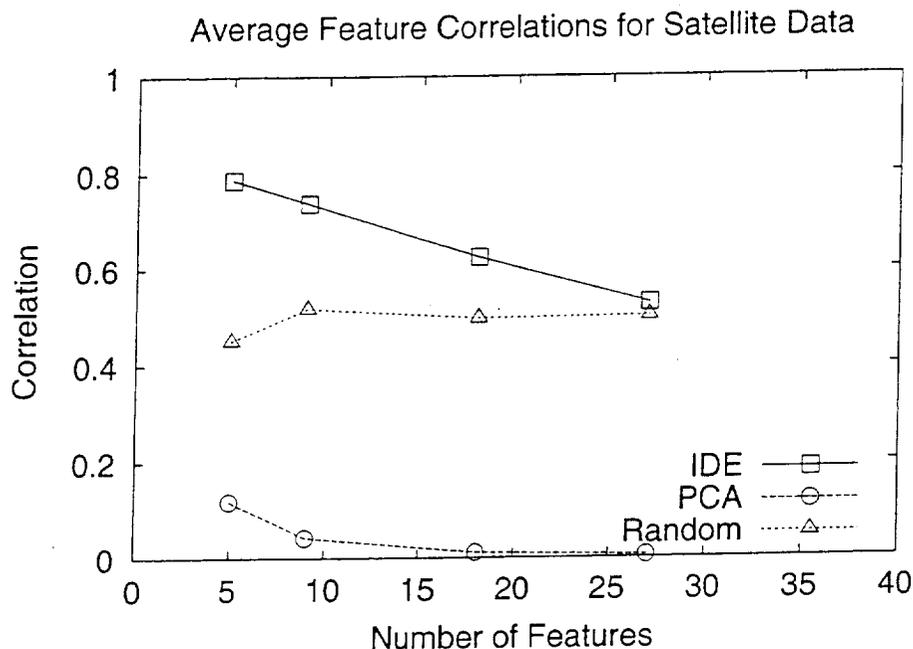


Figure 6: Average Correlation of Features in Satellite Image Data.

three Proben1/UCI datasets, this is the one with the lowest dimensionality, and shows two things: (i) in order to take advantage of input decimation, the initial dimensionality has to be high, as there are likely to be more irrelevant features that can be removed; and (ii) if there are features that have significant meaning, they need to be included in the feature set regardless of their correlation to the particular output. We observed that consecutive groups of four features in the satellite image data set correspond to spectral values for a given pixel. In examining the eigenvalues and eigenvectors, we found that the highest eigenvalue was 91.6% of the sum of the eigenvalues, and the corresponding eigenvector was a simple linear combination of the four spectral values across all the pixels. In this case, the higher principal components provide good discriminative features (i.e., the data “looks” like that in Figure 2(a)). A potential improvement to input decimation is to select “wild card” features based on correlation with all the classes and include them in each decimated subset.

4.3 Synthetic Data

In this section we construct synthetic data sets to enable us to study the properties of input decimated ensembles in a systematic manner. To that end we use the following two synthetic data sets:

- Set A:

- Three classes—one unimodal Gaussian per class.
- 300 training patterns and 150 test patterns—100 training and 50 test patterns per class.
- 100 features per pattern where there are:
 - * 10 relevant features per class. Patterns that belong to a class are generated from a multivariate normal distribution in 10 independent dimensions distributed as $N(40, 5^2)$. There are no dimensions in common among the three classes. Therefore, there are 30 relevant features. For patterns in each class, the 20 features that are relevant to the other two classes are distributed as $U[-100, 100]$.¹⁰
 - * 70 irrelevant features—distributed as $U[-100, 100]$.
- Set B: Same as Set A, except that there is overlap among the relevant features for each class. That is, each class has three relevant features in common with every other class, but there are no features that are relevant to all three classes.

In data set A there is an abundance of features that are irrelevant for the classification task. This data set was chosen to represent large data mining problems where the algorithms may get swamped by irrelevant data. In data set B the overlap among features relevant to each class provides a more difficult problem where the base classifiers are now *forced* to select some common features, reducing the potential for correlation reduction.

4.3.1 Synthetic Set A

Figure 7 presents the classification accuracies and base classifier correlations on Synthetic dataset A as a function of the number of inputs (which are either the number of selected principal components or the number of features selected for each base classifier through input decimation). The original single classifier and original ensemble use all the input features.¹¹ The points for the maximum number of features (e.g., 100 features in this data set), always represent the performance of the original classifier/ensemble.

An important observation that is apparent from these results is that neither PCA ensembles nor PCA base classifiers are particularly sensitive to the number of inputs. The correlations among the base classifiers reinforce this conclusion. Fewer input features in PCA means the base classifiers are more correlated since they all share the same principal features. Note however, that input decimated base classifiers have low correlation for small numbers of features, increasing correlation up to 30 features, and decreasing correlation after that. The base

¹⁰Clearly, because of this, all 30 features have some relevance to all three classes; however, the 10 features used to generate patterns belonging to each class are clearly substantially more relevant than the other 20 features.

¹¹The base classifier used was an MLP with a single hidden layer consisting of 95 units, trained using a learning rate of 0.2 and a momentum term of 0.5.

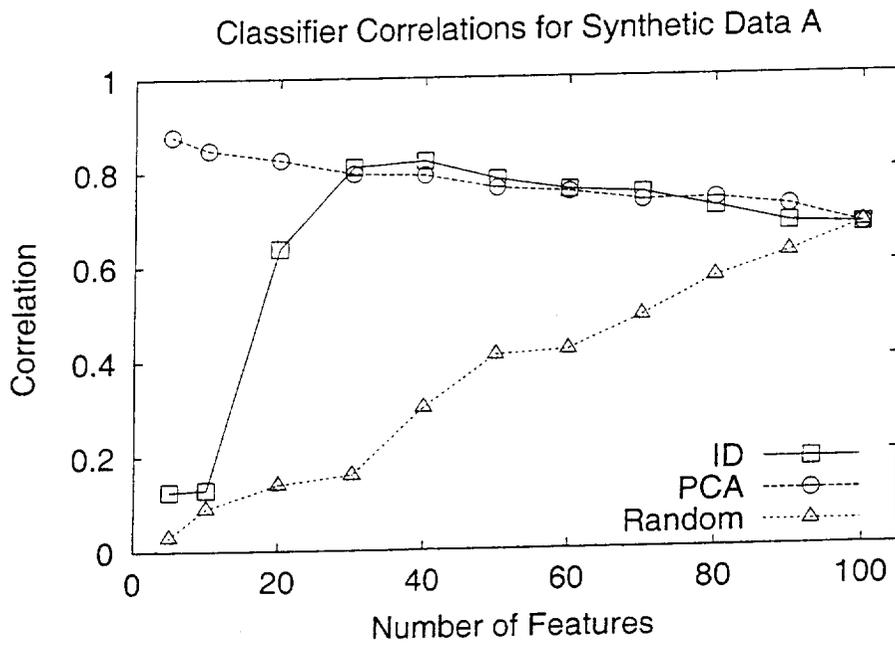
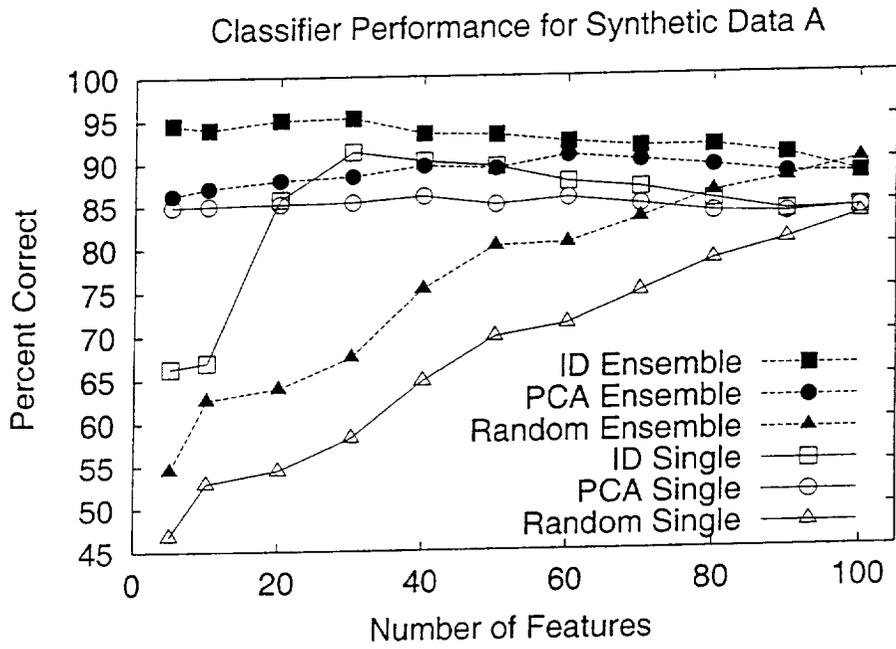


Figure 7: Data set A Performances and Correlations

classifiers' average performance follows a similar pattern. Interestingly, input decimated ensembles are not adversely affected by the poor performance of the base classifiers (e.g., input decimated ensembles with 5 features outperformed input decimated ensembles with 50 features while base classifiers with 5 features gave significantly worse results than base classifiers with 50 features).

In cases where more than 30 features were used, the performance of the ensemble declined with the addition of additional features, i.e., as more and more irrelevant features were included. However, all the input decimation ensembles provided statistically significant improvements over the original ensembles, PCA ensembles, and random-feature ensembles.

The single decimated classifiers with 20 and more features outperformed the original single classifier. This perhaps surprising result (as one might have expected only the ensemble performance to improve when using subsets of the features) is mainly due to the simplification of the learning tasks, which allows the classifiers to learn the mapping more efficiently.

Interestingly, the average correlation among classifiers does not decrease until a very small number of features remain. We attribute this to the removal of noise—removing noise increases the amount of information shared between the base classifiers. Indeed, the correlation increases steadily as features are removed until we reach 30 features (which corresponds to the actual number of relevant features). After that point, removing features reduces the correlation because the base classifiers' feature sets have a decreasing number of common features. The base classifiers' performances also decline; however, the ensemble performance still remains high. This experiment clearly shows a typical trade-off in ensemble learning: one can either increase individual classifier performance (as for input decimation with more than 30 features) or reduce the correlation among classifiers (as for input decimation with less than 20 features) to improve ensemble performance.

4.3.2 Synthetic Set B

Figure 8 presents the results for the second synthetic data set, which is similar to the first data set except that there is overlap among the relevant features for the classes.¹² Because of this overlap, this feature set has fewer total relevant features and thus it constitutes a more difficult problem (as indicated by comparing the results on the full feature base classifiers and ensembles on this data set to the previous one).

Note that the correlations in this data set remained fairly constant across the board for IDE and PCA-based ensembles. Input decimation did not reduce the correlations dramatically for small feature sets in dataset B the way it did in case of dataset A. This is mainly caused by the "coupling" among the base classifiers due to their common input features.

In spite of these difficulties, input decimation ensembles perform extremely well. Indeed, they significantly outperform the original ensemble, PCA ensemble,

¹²The single classifier used was an MLP with a single hidden layer consisting of 95 units, trained using a learning rate of 0.2 and a momentum term of 0.5.

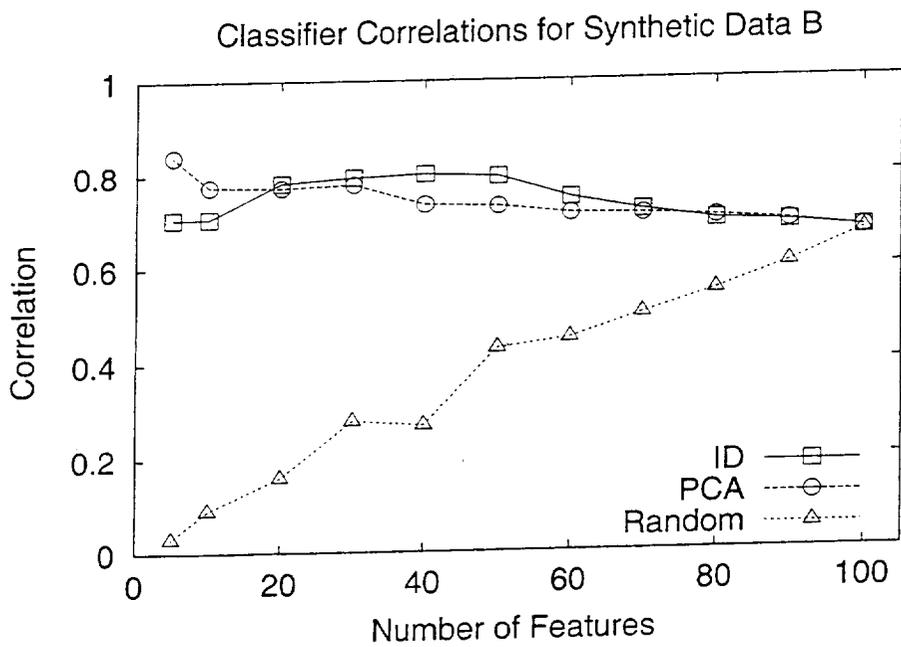
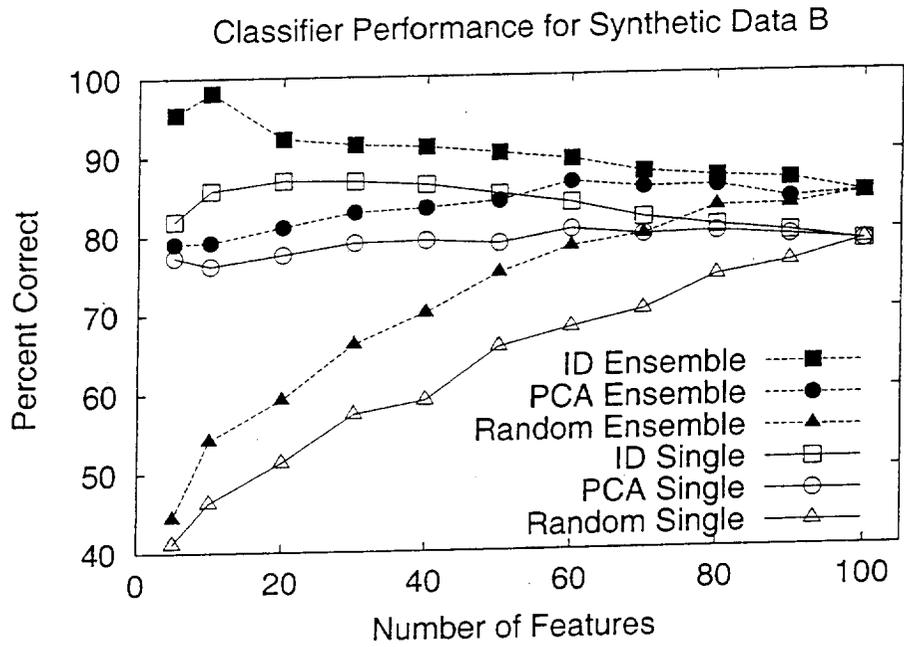


Figure 8: Data set B Performances and Correlations

bles, and random-feature ensembles on all but a few subsets where they only provide marginal improvements. Furthermore the input-decimated single classifiers also outperform their original and PCA counterparts for all but the 60 and 70 feature subsets. This is particularly heartening since this feature set is a more representative abstraction of real data sets (data sets with “clean” separation among classes are quite rare). This experiment demonstrates that when there is overlap among classes, class information becomes particularly relevant. PCA and random feature selection operate without this vital information, therefore they are unlikely to provide competitive performance.

5 Discussion

This paper discusses input decimation, a dimensionality reduction-based ensemble method that provides good generalization performance by reducing the correlations among the classifiers in the ensemble. Through controlled experiments, we show that the input decimated single classifiers often outperform the single original classifiers (trained on the full feature set), demonstrating that simply eliminating irrelevant features can improve performance. In addition, eliminating irrelevant features in each of many classifiers using *different relevance criteria* (in this case, relevance with respect to different classes) yields significant improvement in ensemble performance through correlation reduction, as seen by comparing our decimated ensembles to the original ensembles. Selecting the features using class label information also provides significant performance gains over PCA-based ensembles and random feature subset selection.

Through our tests on synthetic and real data sets, we examined the characteristics that data sets need to have to fully benefit from input decimation. We observed that input decimation yields the greatest improvements over the original ensemble when (i) there are a large number of features (i.e., where it is likely that there will be irrelevant features); and (ii) when the number of training examples is small relative to the input dimensionality (i.e., where it is difficult to properly learn all the parameters in a classifier based on the full feature set). In both cases, by removing the extraneous features, input decimation reduces noise and thereby reduces the number of training examples needed to produce a meaningful model (i.e., alleviating the curse of dimensionality). Our synthetic data sets were generated using multivariate distributions where the feature values were generated independently. We plan to generate synthetic data sets with dependencies among the features to see how they affect our method. Our experiments with real datasets—especially the Satellite Image dataset—showed that input decimation may benefit by keeping out redundant features and including those features that have a high correlation with all classes on average even though they do not have high correlation with any one class. We plan to investigate various possible methods of doing this.

Note that input decimation shares the central aim of generating a diverse pool of classifiers for the ensemble with many methods such as bagging. However, by focusing on the input features rather than the input patterns, input

decimation focuses on a different “axis” of correlation reduction than does bagging. Consequently, input decimation is orthogonal to bagging, and one can use input decimation in conjunction with bagging. We plan to experiment with this in the future.

A final observation is that input decimation works well *in spite* of our rather crude method of feature selection (i.e., using statistical correlation of each feature individually with each class). One reason why this simple method succeeds is that we have greatly simplified the relevance criterion: unlike other feature selection methods that consider the discriminatory ability across all classes, we only consider the relevance of the features to *a single* class. This typically causes each classifier in the ensemble to get a different subset of features, leading to the superior performance we have demonstrated. Nevertheless, we are currently extending this work in four directions: considering cross-correlations among the features; investigating mutual information-based relevance criteria; incorporating global relevance into the selection process; and selecting a different number of features for each classifier.

References

- [1] K. M. Ali and M. J. Pazzani. On the link between error correlation and error reduction in decision tree ensembles. Technical Report 95-38, Department of Information and Computer Science, University of California, Irvine, 1995.
- [2] F.M. Alkoot and J. Kittler. Improving product by moderating k-NN classifiers. In J. Kittler and F. Roli, editors, *Proceedings of the Second International Workshop on Multiple Classifier Systems*, pages 429–439. Springer, Berlin, 2001.
- [3] R. Battiti. Using mutual information for selecting features in supervised neural net learning. *IEEE Transactions on Neural Networks*, 5:4:537–550, July 1994.
- [4] J. O. Berger. *Statistical Decision Theory and Bayesian Analysis*. (2nd Ed.), Springer, New York, 1985.
- [5] C. M. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, New York, 1995.
- [6] C. Blake, E. Keogh, and C.J. Merz. UCI repository of machine learning databases, 1998. (URL: <http://www.ics.uci.edu/~mllearn/MLRepository.html>).
- [7] A. Blum and P. Langley. Selection of relevant features and examples in machine learning. *Artificial Intelligence*, 97:245–272, 1997.

- [8] K. D. Bollacker and J. Ghosh. Linear feature extractors based on mutual information. In *Proceedings of the 13th International Conference on Pattern Recognition*, pages pp. IV:720–724, 1996.
- [9] L. Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996.
- [10] K. J. Cherkauer. Human expert-level performance on a scientific image analysis task by a system using combined artificial neural networks. In *Working Notes of the AAAI Workshop on Integrating Multiple Learned Models*, pages 15–21, 1996.
- [11] S. Cohen and N. Intrator. Automatic model selection in a hybrid perceptron/radial network. In J. Kittler and F. Roli, editors, *Proceedings of the Second International Workshop on Multiple Classifier Systems*, pages 440–454. Springer, Berlin, 2001.
- [12] J.M. Combes, A. Grossman, and Ph. Tchamitchian (Eds.). *Wavelets: Time-Frequency Methods and Phase Space*. Springer-Verlag, 1989.
- [13] D. de Ridder and R. P. W. Duin. Sammon’s mapping using neural networks: A comparison. *Pattern Recognition Letters*, 18:1307–1316, 1997.
- [14] L. Deuser and D. Middleton. On the classification of underwater acoustic signals: An environmentally adaptive approach. *The Acoustic Society of America*, 65:438–443, 1979.
- [15] P.A. Devijver and J. Kittler. *Pattern Recognition: A Statistical Approach*. Prentice-Hall, 1982.
- [16] T.G. Dietterich. Machine learning research: Four current directions. *AI Magazine*, 18(4):97–136, 1998.
- [17] T.G. Dietterich. Ensemble methods in machine learning. In J. Kittler and F. Roli, editors, *Proceedings of the First International Workshop on Multiple Classifier Systems*, pages 1–15. Springer, Berlin, 2000.
- [18] T.G. Dietterich and G. Bakiri. Solving multiclass learning problems via error-correcting output codes. *Journal of Artificial Intelligence Research*, 2:263–286, 1995.
- [19] P. Domingos. Context-sensitive feature selection for lazy learners. *Artificial Intelligence Review*, 11:227–253, 1997.
- [20] R. O. Duda, P. E. Hart, and D. G. Stork. *Pattern Classification*. Wiley, New York, NY, second edition, 2001.
- [21] R.P.W. Duin and D.M.J. Tax. Experiments with classifier combining rules. In J. Kittler and F. Roli, editors, *Proceedings of the First International Workshop on Multiple Classifier Systems*, pages 16–29. Springer, Berlin, 2000.

- [22] Y. Freund and R. Schapire. Experiments with a new boosting algorithm. In *Proceedings of the Thirteenth International Conference on Machine Learning*, pages 148–156. Morgan Kaufmann, 1996.
- [23] K. Fukunaga. *Introduction to Statistical Pattern Recognition*. Academic Press, second edition, 1990.
- [24] J. Ghosh, L. Deuser, and S. Beck. A neural network based hybrid system for detection, characterization and classification of short-duration oceanic signals. *IEEE Journal of Ocean Engineering*, 17(4):351–363, October 1992.
- [25] T. K. Ho. The random space method for constructing decision forests. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(8):832–844, 1998.
- [26] T. K. Ho, J. J. Hull, and S. N. Srihari. Decision combination in multiple classifier systems. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(1):66–76, 1994.
- [27] T.K. Ho. Data complexity analysis for classifier combination. In J. Kittler and F. Roli, editors, *Proceedings of the Second International Workshop on Multiple Classifier Systems*, pages 53–67. Springer, Berlin, 2001.
- [28] A. Hyvarinen. Survey on independent component analysis. *Neural Computing Surveys*, 2:94–128, 1999.
- [29] Robert Jacobs. Method for combining experts' probability assessments. *Neural Computation*, 7(5):867–888, 1995.
- [30] G. H. John, R. Kohavi, and K. Pfleger. Irrelevant features and the subset selection problem. In *Proceedings of the International Conference on Machine Learning (ICML-94)*, pages 121–129, July 1994.
- [31] I.T. Jolliffe. *Principal Component Analysis*. Springer-Verlag, 1986.
- [32] T.M. Jorgensen and C. Linneberg. Feature weighted ensemble classifiers – a modified decision scheme. In J. Kittler and F. Roli, editors, *Proceedings of the Second International Workshop on Multiple Classifier Systems*, pages 218–227. Springer, Berlin, 2001.
- [33] N. Kambhatla and T. K. Leen. Fast non-linear dimension reduction. In J. D. Cowan, G. Tesauro, and J. Alspector, editors, *Advances in Neural Information Processing Systems-6*, pages 152–153. Morgan Kaufmann, 1994.
- [34] N. Kambhatla and T. K. Leen. Dimension reduction by local principal component analysis. *Neural Computation*, 9:1493, 1997.
- [35] J. Kittler. Combining classifiers: A theoretical framework. *Pattern Analysis and Applications*, 1:18–27, 1998.

- [36] J. Kittler and F.M. Alkoot. Relationship of sum and vote fusion strategies. In J. Kittler and F. Roli, editors, *Proceedings of the Second International Workshop on Multiple Classifier Systems*, pages 339–348. Springer, Berlin, 2001.
- [37] J. Kittler, M. Hatef, R.P.W. Duin, and J. Matas. On combining classifiers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(3):226–239, 1998.
- [38] R. Kohavi and G. H. John. Wrappers for feature subset selection. *Artificial Intelligence Journal*, 1-2:273–324, 1997.
- [39] D. Koller and M. Sahami. Toward optimal feature selection. In *Proceedings of the 13th International Conference on Machine Learning*, pages 284–292, 1996.
- [40] A. Krogh and J. Vedelsby. Neural network ensembles, cross validation and active learning. In G. Tesauero, D. S. Touretzky, and T. K. Leen, editors, *Advances in Neural Information Processing Systems-7*, pages 231–238. M.I.T. Press, 1995.
- [41] S. Kumar, M. Crawford, and J. Ghosh. A versatile framework for labelling imagery with a large number of classes. In *Proceedings of the International Joint Conference on Neural Networks (IJCNN-99)*, 1999.
- [42] L.I. Kuncheva and C.J. Whitaker. Feature subsets for classifier combination: An enumerative experiment. In J. Kittler and F. Roli, editors, *Proceedings of the Second International Workshop on Multiple Classifier Systems*, pages 228–237. Springer, Berlin, 2001.
- [43] P. Langley. Selection of relevant features in machine learning. In *Proceedings of the AAAI Fall Symposium on Relevance*, 1994.
- [44] R. Meir. Bias, variance, and the combination of estimators; the case of least linear squares. In G. Tesauero, D. S. Touretzky, and T. K. Leen, editors, *Advances in Neural Information Processing Systems-7*, pages 295–302. M.I.T. Press, 1995.
- [45] C. J. Merz. A principal component approach to combining regression estimates. *Machine Learning*, 36:9–32, 1999.
- [46] M. O. Noordewier, G. G. Towell, and J. W. Shavlik. Training knowledge-based neural networks to recognize genes in DNA sequences. In R.P. Lippmann, J.E. Moody, and D.S. Touretzky, editors, *Advances in Neural Information Processing Systems-3*, pages 530–536. Morgan Kaufmann, 1991.
- [47] E. Oja. *Subspace Methods of Pattern Recognition*. Research Studies Press, Letchworth, England, 1983.

- [48] E. Oja. Principal components, minor components, and linear neural networks. *Neural Networks*, 5:927–936, 1992.
- [49] D. W. Opitz and J. W. Shavlik. Actively searching for an effective neural network ensemble. *Connection Science, Special Issue on Combining Artificial Neural Networks: Ensemble Approaches*, 8(3 & 4):337–354, 1996.
- [50] D. W. Opitz and J. W. Shavlik. Generating accurate and diverse members of a neural-network ensemble. In D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo, editors, *Advances in Neural Information Processing Systems-8*, pages 535–541. M.I.T. Press, 1996.
- [51] N. C. Oza and K. Tumer. Input decimated ensembles: Decorrelation through dimensionality reduction. In J. Kittler and F. Roli, editors, *Proceedings of the Second International Workshop on Multiple Classifier Systems*, pages 238–249. Springer, Berlin, 2001.
- [52] M. Partridge and R. A. Calvo. Fast dimensionality reduction and simple pca. *Intelligent Data Analysis*, 2:203–214, 1998.
- [53] E. Pekalska and R.P.W. Duin. On combining dissimilarity representation. In J. Kittler and F. Roli, editors, *Proceedings of the Second International Workshop on Multiple Classifier Systems*, pages 248–257. Springer, Berlin, 2001.
- [54] Lutz Prechelt. PROBEN1 — A set of benchmarks and benchmarking rules for neural network training algorithms. Technical Report 21/94, Fakultät für Informatik, Universität Karlsruhe, D-76128 Karlsruhe, Germany, September 1994. Anonymous FTP: /pub/papers/tech-reports/1994/1994-21.ps.Z on ftp.ira.uka.de.
- [55] N. Ramanujam, M. F. Mitchell, A. Mahadevan, S. Thomsen, A. Malpica, T. Wright, and N. Atkinson R. Richards-Kortum. Development of a multivariate statistical algorithm to analyze human cervical tissue fluorescence spectra acquired in vivo. *Lasers in Surgery and Medicine*, 19:46–62, 1996.
- [56] M.D. Richard and R.P. Lippmann. Neural network classifiers estimate Bayesian a posteriori probabilities. *Neural Computation*, 3(4):461–483, 1991.
- [57] F. Roli, G. Giacinto, and G. Vernazza. Methods for designing multiple classifier systems. In J. Kittler and F. Roli, editors, *Proceedings of the Second International Workshop on Multiple Classifier Systems*, pages 78–87. Springer, Berlin, 2001.
- [58] B. Rosen. Ensemble learning using decorrelated neural networks. *Connection Science, Special Issue on Combining Artificial Neural Networks: Ensemble Approaches*, 8(3 & 4):373–384, 1996.

- [59] J.W. Sammon Jr. A nonlinear mapping for data structure analysis. *IEEE transactions on Computers*, 18:401–409, 1969.
- [60] N. Short. Remote sensing tutorial, 2000. URL: <http://rst.gsfc.nasa.gov/starthere.html>.
- [61] K. Sirlantzis, M.C. fairhurst, and M.S. Hoque. Genetic algorithms for multi-classifier system configuration: A case study in character recognition. In J. Kittler and F. Roli, editors, *Proceedings of the Second International Workshop on Multiple Classifier Systems*, pages 99–108. Springer, Berlin, 2001.
- [62] D. Tax, M. van Breukelen, R. Duin, and J. Kittler. Combining multiple classifiers by averaging or by multiplying. *Pattern Recognition*, 29(2):341–348, February 1996.
- [63] K. Tumer and J. Ghosh. Analysis of decision boundaries in linearly combined neural classifiers. *Pattern Recognition*, 29(2):341–348, February 1996.
- [64] K. Tumer and J. Ghosh. Error correlation and error reduction in ensemble classifiers. *Connection Science, Special Issue on Combining Artificial Neural Networks: Ensemble Approaches*, 8(3 & 4):385–404, 1996.
- [65] K. Tumer and J. Ghosh. Linear and order statistics combiners for pattern classification. In A. J. C. Sharkey, editor, *Combining Artificial Neural Nets: Ensemble and Modular Multi-Net Systems*, pages 127–162. Springer-Verlag, London, 1999.
- [66] K. Tumer and J. Ghosh. Robust order statistics based ensembles for distributed data mining. In H. Kargupta and P. Chan, editors, *Advances in Distributed and Parallel Knowledge Discovery*, pages 185–210. AAAI/MIT Press, 2000.
- [67] K. Tumer and N. C. Oza. Decimated input ensembles for improved generalization. In *Proceedings of the International Joint Conference on Neural Networks*, 1999.
- [68] D. Windridge and J. Kittler. Classifier combination as a tomographic process. In J. Kittler and F. Roli, editors, *Proceedings of the Second International Workshop on Multiple Classifier Systems*, pages 248–257. Springer, Berlin, 2001.
- [69] D. H. Wolpert. Stacked generalization. *Neural Networks*, 5:241–259, 1992.